## Survey Design & Data Analysis

Identifying Patterns and Drawing Conclusions from Biological Data

_____

_____

_____

_____

_____

_____

_____

_____

## Take Home Lessons

Statistics are used to _simplify_ patterns underlying complex biological phenomena

_Consultation_ with a statistician should be mandatory for any survey-based project

_____

_____

_____

_____

_____

_____

_____

_____

## Take Home Lessons

Statistics are used to _support_ the decision-making process, and not intended to be the sole consideration in that process
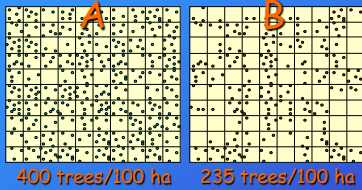
Strong experimental design and statistical analyses lend irrefutable _credibility_ to survey results

_____

_____

_____

_____

_____

_____

_____

_____

## Do we need statistics?

---

### Census

Conclusion:
Site A has a higher density of mahogany than Site B

A          B

400 trees/100 ha    235 trees/100 ha

### Statistics
(using sampling)

Conclusion:
There is a 95% chance that Site A has a higher density than Site B

4.1 ± 0.5 trees/10 ha    2.4 ± 0.6 trees/10 ha

---

## Statistics:

### Descriptive vs Inferential

Describes a sample of the population

Uses sample to generalize to entire population

**Major Steps in designing, Implementing and Evaluating a Project**

1. What is my question or hypothesis?
2. What parameters need to be estimated?
3. Can the parameter be reliably estimated?
4. How will the project be designed?
5. How will the data be analyzed?
6. How will the project be evaluated?

---

**1. What is the Question or Hypothesis?**

*The most important step*

What

Where

When

---

**Sampling Universe**

The population about which you want to draw conclusions

Birds migrating through the Westwoods National Monument

**Question:** Does the abundance of neotropical migratory birds differ among forest interior and young forest/edge habitats at the Westwoods National Monument during Spring migration ?

..................................................

$H_o$: Abundance of NTMBs is the same in both habitat types

$H_A$: Abundance of NTMBs is <u>not</u> the same in both habitat types



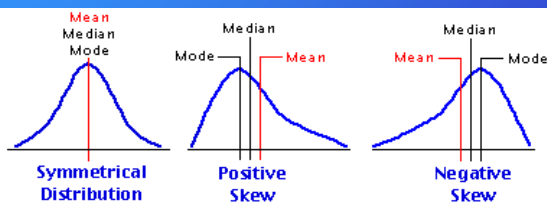## 2. What **Parameter** Needs to be measured?

..................................................

Quantitative characteristic of a population

⇨ Number of individuals

⇨ "Health" of the population

⇨ Environmental threat

⇨ Other characteristics of population

## Types of Biological Data

➡ **Nominal** -- <u>Attribute</u> rather than quantitative
　　　Male, female　　　blue, red, green

➡ **Ordinal** -- Relative difference or <u>ranking</u>
　　　Small, medium, large　　　A, B, C, D,...

➡ **Discrete** -- Quantitative, only <u>whole numbers</u>
　　　0, 1, 2, 3, 4,...

➡ **Continuous** -- Any <u>whole or mixed number</u>
　　　3.4, 19.67, 12.975,...

---

## Ecological data often are <u>not</u> taken from a symmetrical, bell-shaped curve



Symmetrical Distribution　　Positive Skew　　Negative Skew

---

## 3. Can the Parameter be Measured Reliably?

Can you collect enough data to address the question of interest?

Are you measuring the population that you said you would measure?

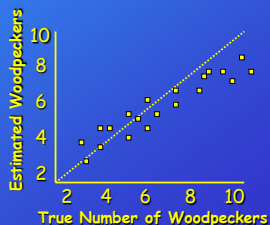Is there excessive ERROR or BIAS in your measurements?

## Sampling

### Error

Random deviations from the true values



### Bias

Systematic deviations from the true values



---

## Major Steps in designing, Implementing and Evaluating a Project

1. What is my question or hypothesis?
2. What parameters need to be estimated?
3. Can the parameter be reliably estimated?
4. How will the project be designed?
5. How will the data be analyzed?
6. How will the project be evaluated?

---

## 4. What is an Appropriate Study Design?

First, go back to Step 2 to review what biological characteristic you are trying to measure

➡ Overall abundance of NTMBs

Second, determine what statistical parameter you need to estimate
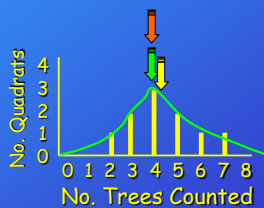
➡ Mean number of NTMBs per plot

## Measures of Central Tendency

*What is most __typical__ for this population*

**Mean** -- average
**Median** -- middle
**Mode** -- most common

No. Quadrats
4
3
2
1
0
0 1 2 3 4 5 6 7 8
No. Trees Counted

**Normal Distribution** -- bell-shaped curve

---

## Measures of Central Tendency

*Most often the MEAN is used
as the parameter of interest*

But, the value of the mean tells
us little without an indication of
the DISPERSION of values used
to calculate that mean.

---

Mean = 4 + 2 + 5 + 4 + 7 + 3 + 5 + 3 + 6 + 4 = 4.3

**Variance ($s^2$) & Standard Deviation (sd)**

$$S^2 = \frac{\Sigma x^2 - \frac{(\Sigma x)^2}{n}}{n-1} = \frac{205 - \frac{(43)^2}{10}}{10-1} = \frac{205 - 184.9}{9}$$

= 20.1 / 9 = 2.23 = Variance

$S^2$ = variance
x = each observation
n = no. of observations
    (sample size)

Standard Deviation =
sd = Sqrt (Variance) =
Sqrt (2.23) = **1.49**

| x | $x^2$ |
|---|---|
| 4 | 16 |
| 2 | 4 |
| 5 | 25 |
| 4 | 16 |
| 7 | 49 |
| 3 | 9 |
| 5 | 25 |
| 3 | 9 |
| 6 | 36 |
| 4 | 16 |

n = 10
$\Sigma x$ = 43
$\Sigma x^2$ = 205

## Standard Deviation (sd)
### What does it mean?

If data follow a normal distribution, then:
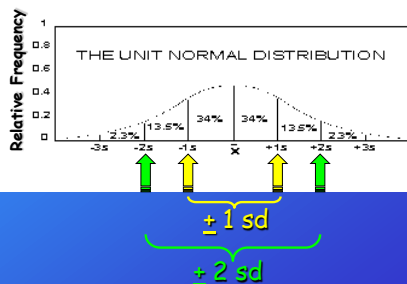
68% of all measurements are within $\pm$ 1 sd

95% of all measurements are within $\pm$ 2 sd

$\overline{x}$ = 4.30 trees/10 ha          sd = 1.49

68% of data are between 4.30 $\pm$ 1.49 (2.81 - 5.79)

95% of data are between 4.30 $\pm$ 2.98 (1.32 - 7.28)

---

### Probability and the Normal Distribution

THE UNIT NORMAL DISTRIBUTION

Relative Frequency

2.3%   13.5%   34%   34%   13.5%   2.3%

-3s   -2s   -1s   $\overline{x}$   +1s   +2s   +3s

$\pm$ 1 sd

$\pm$ 2 sd

---

## Standard Deviation (sd)
### vs
## Standard Error (SE)

sd = variation in the population (sample)

SE = how close estimated mean
is to the true population mean

$$SE = sd / \sqrt{n}$$

The Standard Deviation is good for examining variation around the mean, but what if we want to compare the variation in two populations that differ widely in their mean values?



**Does one species show greater variation in weight?**

Asian Elephant
x= 3960 kg
sd = 283

Elephant Shrew
x = 0.22 kg
sd = 0.10

## Coefficient of Variation CV)

Provides a measure of relative variability

$$CV = (sd / \bar{x})(100\%)$$



(283/3960)(100%) = 7%



(0.10/0.22)(100%) = 5%

## Standard Deviation (sd)
## vs
## Standard Error (SE)
## vs
## Coefficient of Variation (CV)

sd = variation in the population (sample)

SE = how close estimated mean is to the true population mean = sd / sqrt(n)

CV = relative estimate of population variability = sd / mean

## 4. What is an Appropriate Study Design?

✔ First, go back to Step 2 to review what biological characteristic you are trying to measure

✔ Second, determine what statistical parameter you need to estimate

Third, develop sampling protocol that allows reliable data to be collected

What will be our sampling protocol?

Line transects?

Mist netting?

Spot mapping?

Point counts?

## Sampling Design

Randomness -- Every unit in the population has an equal chance of being sampled.

Independence -- Knowing something about one unit doesn't provide information about another unit (or one unit does not influence another unit).
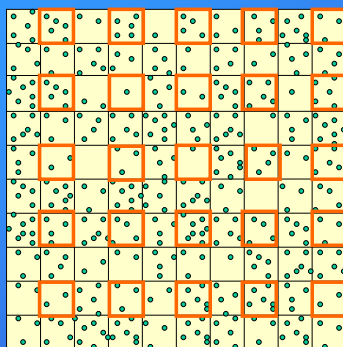
## Random (or Probability) Sampling



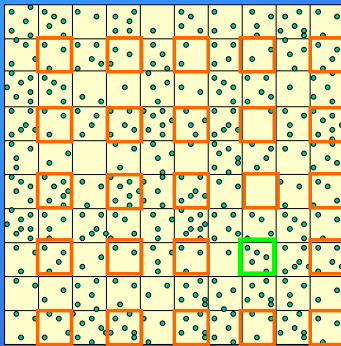Quadrats selected totally by chance

Random numbers table

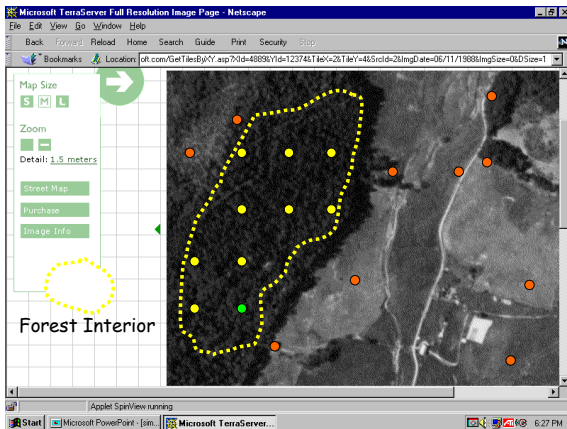## Systematic Sampling



Quadrats selected are evenly spaced

No random component

Systematic Sampling with Random Start

Quadrats selected are evenly spaced

Random component



Forest Interior

Minimum Sample Size Requirements

or

How do we know
how many point count stations
we might need?

## Minimum Sample Size Requirements

Calculated BEFORE study is implemented
(IF POSSIBLE)

Use PILOT STUDY or existing data

Need 4 pieces of information:
1. Mean & measure of variance of parameter
2. Magnitude of difference you want to detect
3. Significance level ($\alpha$)
4. Desired power

---

## Some Statistical Terminology

�֎ Type I statistical error

✖ Type II statistical error

✖ Power

---

## Statistical Errors in Hypothesis Testing

Alpha ($\alpha$) -- Also called Type I error.
Probability that we reject the
Null Hypothesis when in fact it is true.

Beta ($\beta$) -- Also called Type II error.
Probability that we do not reject the
Null Hypothesis when it is, in fact, false.

## The **Power** of Statistical Tests

The probability of rejecting the
Null Hypothesis when, in fact,
it is false (and should be rejected)

Power = $1 - \beta$

Important for calculating minimum
necessary sample sizes

## Minimum Sample Size Requirements

Calculated BEFORE study is implemented
(IF POSSIBLE)

Use PILOT STUDY or existing data

Need 4 pieces of information:
1. Mean & measure of variance of parameter
2. Magnitude of difference you want to detect
3. Significance level ($\alpha$)
4. Desired power

## Minimum Sample Size (n)

$$n = \frac{2Ms^2}{d^2}$$

Where, d = minimum detectable difference

M = multiplier from normal distribution

$s^2$ = estimated population variance

# 5. Appropriate Data Analysis

Good analyses begin
with good hypotheses

All statistical tests must have
two types of hypotheses:

Null Hypothesis ($H_o$)
Alternative Hypothesis ($H_A$)

---

Null Hypothesis usually is
tested via a statistical test

If Null Hypothesis not accepted,
then Alternative Hypothesis
is assumed to be true

We need an objective way of rejecting
or not rejecting the null hypothesis based
upon the probability that the estimated
parameter occurred by chance alone.

---

The conclusion that the observed
result is significant is established
by the significance level (alpha or $\alpha$).

It is the probability above which we
do not reject the Null Hypothesis.

Now, let's get down to business...

_____

_____

_____

_____

_____

_____

_____


**Question**: Does the abundance of neotropical migratory birds differ among forest interior and young forest/edge habitats at the Westwoods National Monument during Spring migration?

...........................................................

$H_o$: Abundance of NTMBs is the same in both habitat types ($\alpha$ = 0.05)

$H_A$: Abundance of NTMBs is <u>not</u> the same in both habitat types

_____

_____

_____

_____

_____

_____

_____

_____


Let's take a quick look at the data...

| Plot | Forest Interior | | Young Forest/Edge | |
|------|-----------------|---|-------------------|---|
| 1 | 3 | | 4 | |
| 2 | 7 | | 5 | |
| 3 | 4 | | 5 | |
| 4 | 3 | | 5 | |
| 5 | 2 | | 4 | |
| 6 | 1 | | 2 | |
| 7 | 4 | | 3 | |
| 8 | 5 | | 6 | |
| 9 | 3 | $\bar{x}$ = 3.60 | 6 | $\bar{x}$ = 4.90 |
| 10 | 4 | $s$ = 1.65 | 9 | $s$ = 1.91 |

_____

_____

_____

_____

_____

_____

_____

_____

### FOREST INTERIOR
Mean = 3 + 7 + 4 + 3 + 2 + 1 + 4 + 5 + 3 + 4 = 3.60

**Variance ($s^2$) & Standard Deviation (sd)**

$$S^2 = \frac{\Sigma x^2 - \frac{(\Sigma x)^2}{n}}{n-1} = 2.72 = \text{Variance}$$

$S^2$ = variance
x = each observation
n = no. of observations
(sample size)

Standard Deviation =
sd = Sqrt (Variance) =
Sqrt (2.72) = **1.65**

| x | $x^2$ |
|---|-------|
| 3 | 9 |
| 7 | 49 |
| 4 | 16 |
| 3 | 9 |
| 2 | 4 |
| 1 | 1 |
| 4 | 16 |
| 5 | 25 |
| 3 | 9 |
| 4 | 16 |

n = 10
$\Sigma x$ = 36
$\Sigma x^2$ = 154

---

## Probability and the Normal Distribution



THE UNIT NORMAL DISTRIBUTION

2.3%  13.5%  34%  34%  13.5%  2.3%

-3s  -2s  -1s  x̄  +1s  +2s  +3s

± 1 sd ---------> 68%

± 2 sd ---------> 95%

---

## Are the Data Normally Distributed?

**Forest Interior**

x̄ = 3.60
s = 1.65



**Young Forest /Edge**

x̄ = 4.90
s = 1.91

So…
The data appear to follow an
(approximate) <u>normal distribution</u>

So…
We can use <u>parametric</u> statistics
to analyze the data

_____

We choose a <u>t-test</u>, because:

We are comparing only <u>two</u> samples

The data are <u>normally distributed</u>

The two samples have <u>similar variances</u>

_____

The <u>t-test</u> will allow us to assess
if there is a difference in the abundance
of NTMBs in the two habitat types,
Forest Interior and Young Forest/Edge

$$t = \dfrac{\overline{X} - \overline{X}}{\sqrt{\dfrac{\left(\dfrac{\Sigma x^2 - \dfrac{(\Sigma x)^2}{n}}{n-1}\right)}{n} + \dfrac{\left(\dfrac{\Sigma x^2 - \dfrac{(\Sigma x)^2}{n}}{n-1}\right)}{n}}}$$

## The test statistic...

t = 1.63



Look up critical value in table: $\alpha$ = 0.05 (2-tailed), $v$ = 18

_____

_____

_____

_____

_____

_____

_____

_____

_____

## What do we conclude?

$H_o$: Abundance of NTMBs is the same in both habitat types

$H_A$: Abundance of NTMBs is not the same in both habitat types
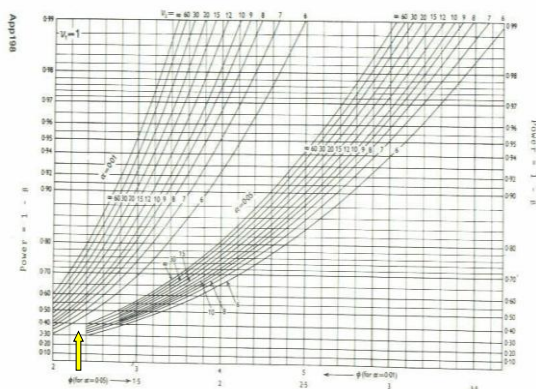
Since the t-statistic (1.63) is not greater than the critical value from the table (2.101), we cannot reject the null hypothesis and conclude that no difference exists between habitats.

_____

_____

_____

_____

_____

_____

_____

_____

**Another way to examine the differences between the mean values of two samples**

_____

_____

_____

_____

_____

_____

_____

_____

_____

## Confidence Intervals

An estimated range of values which is likely to include an unknown population parameter

Often, 95%

Allows us to estimate the precision of our estimate

---

*t* value taken from statistical table

$$\overline{x} \pm (t_{\alpha(2),n-1}) \left( \frac{sd}{\sqrt{n}} \right)$$

Estimated mean

Standard deviation divided by square root of sample size

---

### Mean + 95% Confidence Intervals



Birds per Station

Interior    Edge

## Now, back to our hypothess testing

What do we conclude?

$H_o$: Abundance of NTMBs is the same in both habitat types

$H_A$: Abundance of NTMBs is not the same in both habitat types

Since the t-statistic (1.63) is not greater than the critical value form the table (2.101), we cannot reject the null hypothesis and conclude that no difference exists between habitats.

---

## What do we conclude?

$H_o$: Abundance of NTMBs is the same in both habitat types

$H_A$: Abundance of NTMBs is not the same in both habitat types

Since the t-statistic (1.63) is not greater than the critical value form the table (2.101), we cannot reject the null hypothesis and conclude that no difference exists between habitats.

---

But, was our design rigorous enough to be able to detect a difference?

In other words, what was our Power and Type II error?

$$\phi = \sqrt{\frac{n \, (\text{difference}^2)}{4 \, (\text{variance})}}$$

= 1.15    Look up Power in Figure B.1a

The Power of this statistical test was approximately 0.30, which means...

That there was only a 30% chance that we could have detected this difference

or

that there was a 70% chance that we claimed there was no difference in NTMB abundance when, in fact, there was a difference

How many point counts would we have needed in each habitat to detect a 26% difference?

$\dfrac{4.90 - 3.60}{4.90}$   $n = \dfrac{2Ms^2}{d^2}$

M = multiplier from table = 7.9
$s^2$ = variance = 3.648
d = difference = 4.90 - 3.60 = 1.30

$n = \dfrac{2\,(7.9)\,(3.648)}{(1.30)(1.30)} = \dfrac{57.64}{1.69} = 34.1$ points

Browser window: **Power Analysis of Monitoring Programs - Netscape** http://im.nbs.gov/powcase/powcase.html

Location: http://www.im.nbs.gov/powcase/powcase.html

*Power Analysis of Monitoring Programs*
**Designing effective surveys**

Main Page
A Power Primer
Power Analysis Step-by-Step
How many, how often, & how long?
Power Links
Power Bibliography
Download MONITOR
Thanks to...

Inherent to any monitoring program is variation in the numbers of plants and animals counted. Some of that variation is natural (e.g., population dynamics such as births, deaths, immigration, and emigration, weather effects, and so on) and some is due to the flaws of the chosen monitoring technique (e.g., observer differences, different fractions of individuals being counted each time). This variation in numbers (from both natural and sampling sources) partially obscures the presence of any long-term trends. If the noise from this extraneous variation is high enough and we have too few counts, then we may fail to detect important underlying population trends in our population which, after all, is the goal of our monitoring program.

These pages are intended for anyone who's interested in starting a monitoring program, and do not assume advanced knowledge of statistics.

Other sites of interest in the *PWRC Monitoring Program* family:
The Amphibian Count Database

---

## Conclusions

No significant difference existed in abundance of NTMBs between forest interior and young forest/edge habitats...

but

we had only a 30% chance of detecting a difference if it did, indeed, exist

---

## 6.  Evaluating Success of Project

Completeness of data

Accuracy of data

Appropriateness of data

"Adaptive" Approach

## Data Management

All survey data should be:

Checked for errors before & after recording in an electronic format

Recorded in an electronic format ASAP

Stored in two separate locations

Accompanied by metadata

---

## Choosing an Appropriate Statistical Test

Type of Data

Number of Variables

Sample Characteristics

Nature of hypothesis/research question

---

### How many variables?

One Variable

One Group — Two Groups — 3+ Groups

Normal?   yes      no         yes        no         yes        no

mean & SD    binomial    t-test    Chi-square    ANOVA    non-parametric

## How many variables?



## How many variables?



## Example

# Lizards…

## …of the Venezuelan Savanna

**We are interested in detecting:**

A change in lizard population density of at least 50%

At a significance ($\alpha$) level of 0.10

And power of 90%

**Minimum Sample Size (n)**

$$n = \frac{2Ms^2}{d^2}$$

Where, d = minimum detectable difference

M = multiplier from normal distribution

$s^2$ = estimated population variance



**Minimum sample size,** $n = \dfrac{2Ms^2}{d^2}$

M = 8.6    Mean = 0.60    $s^2$ = 0.49

d = (0.50)(0.60) = 0.30

$$n = \frac{2(8.6)(0.49)}{(0.30)^2} = \frac{8.43}{0.09} = 94$$

**CORRECTED Minimum Sample Size (n')**

Because we plan
to sample such
a large proportion
of total area

---

**CORRECTED Minimum Sample Size (n')**

$$n' = \frac{n}{(1 + [n/N])}$$

n' = corrected sample size

n = original sample size

N = total possible sample size

$$n' = \frac{94}{(1 + [94/100])} = \frac{94}{1.94} = \mathbf{48}$$

---



100 m
100 m

Now, let's select our plots

**We're also interested in habitat use by lizards in Llanos National Park**

**Question: Do lizards exhibit a habitat preference in Llanos National Park?**

**$H_o$: No difference in lizard abundance between habitats**

**$H_A$: Lizards not distributed equally between habitats**

Alpha = 0.05

---

**Use Mann-Whitney U-Test to test the null hypothesis**

**How to calculate U**

**Step 1. Rank order all counts of lizards from each of the two habitats**

**Step 2. Sum the ranks from the smaller sample. This gives $R_1$. [595.5]**

---

**Step 3. Calculate $U_1$ from the equation:**

$$U_1 = \frac{(n_1)(n_2) + n_1(n_1 + 1)}{2} - R_1$$

**Where,**
**$n_1$ = sample size for sample 1 [21]**
**$n_2$ = sample size for sample 2 [27]**

**$U_1 = 202.5$**

**Step 4.** Calculate $U_2$ from the equation:

$$U_2 = (n_1)(n_2) - U_1 \quad [U_2 = 364.5]$$

**Step 5.** Take the larger of $U_1$ & $U_2$ and call that U. With small sample sizes, you can compare U to values in a statistical table. But, with large sample sizes, the hypothesis must be tested using a normal approximation.



Two-tailed test, P = 0.05

One-tailed test, P = 0.05

0.05

0.025

0.025

−1.96     0     1.64  1.96     z-score

So, let's calculate our Z score
And see where it falls along the
normal distribution curve

$$Z = \frac{U - \mu_U}{\sigma_U} \text{ , where} \quad [Z = 1.68]$$

$$\mu_U = \frac{(n_1)(n_2)}{2} \qquad \sigma_U = \sqrt{\frac{(n_1)(n_2)(N+1)}{12}}$$

Appendix B    Statistical Tables and Graphs    App17

TABLE B.2   Proportions of the Normal Curve (One-Tailed)

**Z score must be ≥ 1.96**

**Observed z-score**



**One-tailed test, P = 0.05**

**Two-tailed test, P = 0.05**

0.05

0.025            0.025

−1.96        0        1.64   1.96    z-score



**What do we conclude about the abundance of lizards in the two habitat types?**

**How might you better design this study?**